# On Verifying and Engineering the Well-gradedness of a Union-closed Family[*]

David Eppstein[†]     Jean-Claude Falmagne[‡]

{eppstein, jcf}@uci.edu

Hasan Uzun[§]

huzun@aleks.com

April 14, 2008

## Abstract

Current techniques for generating a knowledge space, such as QUERY, guarantees that the resulting structure is closed under union, but not that it satisfies wellgradedness, which is one of the defining conditions for a learning space. We give necessary and sufficient conditions on the base of a union-closed set family that ensures that the family is well-graded. We consider two cases, depending on whether or not the family contains the empty set. We also provide algorithms for efficiently testing these conditions, and for augmenting a set family in a minimal way to one that satisfies these conditions.

## Introduction

A family of sets $\mathcal{F}$ is well-graded if any two sets in $\mathcal{F}$ can be connected by a sequence of sets formed by single-element insertions and deletions, without redundant operations, such that all intermediate sets in the sequence belong to $\mathcal{F}$. The family $\mathcal{F}$ is called ∪-closed if it is closed under union. (Formal definitions are given in our next section.) Well-graded families are of interest for theorists in several different areas of combinatorics, as various families of sets or relations are well-graded. For example, Theorems 2 and 4 in Bogart (1973) imply that the family of all partial orders on a finite set is well-graded. The same property of well-gradedness is shared by other families, such as the semiorders, the interval orders, and the biorders, again on finite sets (Doignon and Falmagne, 1997). Via representation theorems, this concept also applies to the *partial cubes*, to wit, graphs isometrically embeddable into hypercubes (Graham and Pollak, 1971; Djoković, 1973; Winkler, 1984; Imrich and Klavžar, 2000),

---

[†]Computer Science Department, University of California, Irvine, CA 92697.

[‡]Dept. of Cognitive Sciences, University of California, Irvine, CA 92697. Phone: (949) 433 2735.

[§]ALEKS Corporation.

and to the *oriented media* which are semigroups of transformations satisfying certain axioms (Falmagne, 1997; Eppstein et al., 2007).

When the family $\mathcal{F}$ is well-graded, $\cup$-closed and contains the empty set, one obtains an object variously called an *antimatroid* (Korte et al., 1991), a *learning space* (Cosyn and Uzun, 2008; Falmagne et al., 2006), or a *well-graded knowledge space* (Doignon and Falmagne, 1985). The monograph of Doignon and Falmagne (1999) contains a comprehensive account of this topic. Learning spaces are applied in mathematical modeling of education. In such cases, the ground set is the collection of problems, for example in elementary arithmetic, that a student must learn to solve in order to master the subject. The family $\mathcal{F}$ contains then all the subsets forming the feasible *knowledge states*. In practice, the size of such a family is quite large, typically containing millions of states[1], which raises the problem of summarizing $\mathcal{F}$ efficiently. An obvious choice for this purpose is the *base* of that family, namely the unique minimal subset of $\mathcal{F}$ whose completion via all possible unions gives back $\mathcal{F}$.

For various reasons, when building a learning space in practice, one may fall short of some sets to achieve well-gradedness, a property regarded as essential for promoting efficient learning (see the axiomatization of Cosyn and Uzun, 2008). This raises the problems of uncovering possibly missing sets, and completing the family economically and/or optimally. These considerations inspired the work presented here.

We solve the following problems for a finite family $\mathcal{G}$ of finite sets.

1. Find necessary and sufficient conditions for $\mathcal{G}$ to be the base of a well-graded $\cup$-closed family of sets.

2. Find such conditions when the well-graded $\cup$-closed family of sets is known to be a learning space, that is, the family contains the empty set. (These conditions may be simpler than in Case 1.)

3. Provide efficient algorithms for testing these conditions on a family $\mathcal{G}$ and uncovering possibly missing sets. (Different algorithms may be used in Problems 1 and 2.)

4. Supposing that some family $\mathcal{G}$ fails to satisfy the conditions in Problems 1 or 2, provide algorithms for modifying $\mathcal{G}$ in some optimal sense to yield a family $\mathcal{G}^*$ satisfying such conditions.

Except for the passing remark involving Counterexample 12, only finite sets are considered in this paper.

# Background and Preparatory Results

**1 Definition.** Let $\mathcal{F}$ be a family of subsets of a set $\mathcal{X}$. A *tight path between* two distinct sets $P$ and $Q$ (or from $P$ to $Q$) in $\mathcal{F}$ is a sequence $P_0 = P, P_1, \ldots, P_n = Q$ in $\mathcal{F}$ such that $d(P,Q) = |P \triangle Q| = n$ and $d(P_i, P_{i+1}) = 1$ for $0 \leq i \leq n-1$.

---

[1] For a ground set that may contain a couple of hundreds of problem types.

The family $\mathcal{F}$ is *well-graded* or a *wg-family* if there is a tight path between any two of its distinct sets. (cf. Doignon and Falmagne, 1997; Falmagne and Doignon, 1997)[2].

**2 Definition.** A family of sets $\mathcal{F}$ is *closed under union*, or $\cup$-*closed* if for any nonempty[3] $\mathcal{G} \subseteq \mathcal{F}$ we have $\cup\mathcal{G} \in \mathcal{F}$. A well-graded family closed under union and containing the empty set is a *learning space*.

**3 Definition.** The *span* of a family of sets $\mathcal{G}$ is the family $\mathcal{G}^\dagger$ containing any set which is the union of some subfamily[4] of $\mathcal{G}$. In such a case, we write $\mathbb{S}(\mathcal{G}) = \mathcal{G}^\dagger$ and we say that $\mathcal{G}$ *spans* $\mathcal{G}^\dagger$. By definition $\mathbb{S}(\mathcal{G})$ is thus $\cup$-closed. A *base* of a $\cup$-closed family $\mathcal{F}$ is a minimal subfamily $\mathcal{B}$ of $\mathcal{F}$ spanning $\mathcal{F}$ (where 'minimal' is meant with respect to set inclusion: if $\mathbb{S}(\mathcal{H}) = \mathcal{F}$ for some $\mathcal{H} \subseteq \mathcal{B}$, then $\mathcal{H} = \mathcal{B}$). Notice that if $\varnothing \in \mathcal{F}$, we must have $\varnothing \in \mathcal{B}$, with $\cup\{\varnothing\} = \varnothing$. In such a case, we use the abbreviation $\breve{\mathcal{B}} = \mathcal{B} \setminus \{\varnothing\}$. Note that a family $\mathcal{G}$ spanning a family $\mathcal{F}$ is a base of $\mathcal{F}$ if and only if none of the sets in $\mathcal{G}$ is the union of some other sets in $\mathcal{G}$.

Any finite $\cup$-closed family has a base, which is unique. This uniqueness property of the base also holds in the infinite case but some infinite families have no base: take, for example, the collection of all open sets of $\mathbb{R}$ or Counterexample 12. (In this regard, see Doignon and Falmagne, 1999, Theorems 1.20 and 1.22.)

The following lemma is a key tool, as it allows us to infer the wellgradedness of a family from that of its base.

**4 Lemma.** *The span of a wg-family is well-graded.*

PROOF. Let $\mathbb{S}(\mathcal{G})$ be the span of some wg-family $\mathcal{G}$. Take any two distinct $X, Y$ in $\mathbb{S}(\mathcal{G})$. Since $\mathbb{S}(\mathcal{G})$ is $\cup$-closed by definition, $X \cup Y$ is in $\mathbb{S}(\mathcal{G})$ and we have $d(X, Y) = d(X, X \cup Y) + d(X \cup Y, Y)$. Accordingly, it suffices to prove that there is in $\mathbb{S}(\mathcal{G})$ a tight path

$$(1) \qquad X_1 = X, X_2, \ldots, X_n = X \cup Y,$$

with in fact $X_i \subset X_{i+1}$, $1 \leq i \leq n - 1$. By definition of the span, there exists finite $\mathcal{H}, \mathcal{K} \subseteq \mathcal{G}$ such that $X = \cup\mathcal{H}$ and $Y = \cup\mathcal{K}$. Without loss of generality (exchanging the roles of $X$ and $Y$ if needed), we can assume that there exists some $K \in \mathcal{K}$ such that $K \setminus X \neq \varnothing$. Choose $H \in \mathcal{H}$ arbitrarily. By the wellgradedness of $\mathcal{G}$, there is a tight path $H_1 = H, \ldots, H_m = K$. Let $k$ be the first index such that $H_k \setminus X \neq \varnothing$. (Such an index must exist because $K \setminus X \neq \varnothing$.) We necessarily have $|H_k \setminus X| = 1$. Defining $X_2 = (\cup\mathcal{H}) \cup H_k$, we obtain $X_1 = X \subset X_2 \subseteq X \cup Y$ with $|X_2 \setminus X_1| = 1$. An induction completes the proof. □

Note however that the base of a $\cup$-closed wg-family need not be well-graded.

---

[2]This concept was introduced earlier under a different name; see Kuzmin and Ovchinnikov (1975), Ovchinnikov (1980).

[3]For some authors, the subfamily $\mathcal{G}$ may be empty, with $\cup\varnothing = \varnothing$. So, a $\cup$-closed family automatically contains the empty set. We do not use this convention here.

[4]Contrary to the convention used by Doignon and Falmagne (1999), the empty subfamily of $\mathcal{G}$ is not allowed; so $\varnothing \in \mathbb{S}(\mathcal{G})$ only if $\varnothing \in \mathcal{G}$.

**5 Example.** The $\cup$-closed wg-family

$$\mathcal{F} = \{\varnothing, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \{c, d\}, \{a, b, c\},$$
(2)
$$\{a, c, d\}, \{b, c, d\}, \{a, b, c, d\}, \{a, b, c, d, e\}\}.$$

has the base $\{\varnothing, \{a\}, \{b\}, \{c\}, \{c, d\}, \{a, b, c, d, e\}\}$, which is not well-graded. Moreover, $\mathcal{F}$ has two different minimal well-graded subfamilies spanning $\mathcal{F}$:

$$\{\varnothing, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{c, d\}, \{a, b, c\},$$
(3)
$$\{a, c, d\}, \{a, b, c, d\}, \{a, b, c, d, e\}\},$$
$$\{\varnothing, \{a\}, \{b\}, \{c\}, \{a, b\}, \{b, c\}, \{c, d\}, \{a, b, c\},$$
(4)
$$\{b, c, d\}, \{a, b, c, d\}, \{a, b, c, d, e\}\}.$$

**6 Example.** Notice that the base of a family which is both $\cup$-closed and $\cap$-closed (that is, closed under intersection) is not necessarily well-graded. Indeed, consider the family

$$\mathcal{G} = \{\varnothing, \{a\}, \{b\}, \{d\}, \{a, b\}, \{a, d\}, \{b, d\}, \{a, b, c\}, \{a, b, d\},$$
$$\{a, b, c, d\}, \{a, b, c, d, e\}\},$$

for which $\{\varnothing, \{a\}, \{b\}, \{d\}, \{a, b, c\}, \{a, b, c, d, e\}\}$ is the base.

# Main Results

**7 Theorem.** *Let $\mathcal{F}$ be a $\cup$-closed family with base $\mathcal{B}$. Then $\mathcal{F}$ is a wg-family if and only if, for any two distinct sets $K$ and $L$ in $\mathcal{B}$, there is a tight path in $\mathcal{F}$ from $K$ to $L \cup K$. If $\mathcal{B}$ contains the empty set, then $\mathcal{F}$ is well-graded if and only if there is a tight path from $\varnothing$ to $K$ for any $K$ in $\mathcal{B}$.*

Thus, this result provides a solution to Problems 1 and 2. Another solution to Problem 2 is given by Lemma 19.

PROOF. As $\mathcal{F}$ is $\cup$-closed with base $\mathcal{B}$, the necessity is clear for both statements. To establish that the sufficiency in the first statement also holds, we point out that the family $\mathcal{B}^*$ defined by

(5)
$$M \in \mathcal{B}^* \iff \begin{cases} M = \cup \mathcal{A} \text{ for some } \mathcal{A} \subseteq \mathcal{B} \text{ such that} \\ K \subseteq \cup \mathcal{A} \subseteq K \cup L \text{ for some } K, L \in \mathcal{B} \end{cases}$$

includes $\mathcal{B}$ since $K = \cup\{K\}$ and $K \subseteq \cup\{K\} \subseteq K \cup L$ for any $K$ and $L$ in $\mathcal{B}$. Since $\mathcal{B} \subseteq \mathcal{B}^* \subseteq \mathcal{F}$ the family $\mathcal{B}^*$ spans $\mathcal{F}$. We claim that $\mathcal{B}^*$ is well-graded, which implies by Lemma 4 that $\mathcal{F}$ is well-graded. The main line of our argument is similar to that used in the proof of Lemma 4.

Take any two distinct $V, W \in \mathcal{B}^*$. By definition of $\mathcal{B}^*$, we have $V = \cup \mathcal{V}$ and $W = \cup \mathcal{W}$ for some subfamilies $\mathcal{V}$ and $\mathcal{W}$ of $\mathcal{B}$. Suppose that $d(V, V \cup W) = n$. We have to show that there exists in $\mathcal{B}^*$ a tight path

$$V_0 = V, V_1, \ldots, V_n = V \cup W$$

from $V$ to $V \cup W$. Without loss of generality (exchanging the roles of $V$ and $W$ if needed), we can assume that there is some $H \in \mathcal{W}$ such that $H \setminus V \neq \varnothing$. Choose $G \in \mathcal{V}$ arbitrarily. Then $G \subset H \cup G \subseteq V \cup W$, with $H$ and $G$ in $\mathcal{B}$. By hypothesis, there is a tight path $G_0 = G, G_1, \ldots, G_m = G \cup H$ from $G$ to $G \cup H$ in $\mathcal{F}$, with $G \subset G_i \subset G \cup H$ and $d(G, G_i) = i$ for $1 \leq i \leq m$. Let $k$ be the first index such that $G_k \setminus V \neq \varnothing$. (Such an index must exist because $H \setminus V \neq \varnothing$.) We necessarily have $|G_k \setminus V| = 1$. Defining $V_1 = (\cup \mathcal{V}) \cup G_k$, we obtain $V_0 = V \subset V_1 \subseteq V \cup W$ with $|V_1 \setminus V_0| = 1$. An induction completes the proof of the sufficiency for the first statement.

We now show that if $\varnothing \in \mathcal{B}$, then there is a tight path from $L$ to $K \cup L$ for any $K$ and $L$ in $\mathcal{B}$. Thus, the sufficiency of the second statement follows from that in the first statement. Indeed, let $K_0 = \varnothing, K_1, \ldots, K_n = K$ be a tight path. It is easily seen that, after removal of identical terms if need be, the sequence $K_0 \cup L = L$, $K_1 \cup L, \ldots, K_n \cup L = K \cup L$ is a tight path from $L$ to $K \cup L$. $\qquad \square$

**8 Remark.** The set $\mathcal{B}^*$ constructed in the proof of Theorem 7 is not necessarily a minimal wg-family spanning $\mathcal{F}$. Indeed, the definition of $\mathcal{B}^*$ by (5) includes all the unions $\cup \mathcal{A}$, while only some of them may be needed. An example was provided by the wg-family of Example 5. In this case, each of (3) and (4) is a minimal wg-family including the base and spanning the wg-family $\mathcal{F}$ defined by (2). The set $\mathcal{B}^*$ in this case would be the union of the two families in (3) and (4), which is in fact equal to $\mathcal{F}$.

In the case of learning spaces, Koppen (1998) obtained a different, but equivalent answer to Problem 1 (see Theorem 13). As shown by Counterexample 14, Koppen's result does not generalize to the case in which the family does not contain the empty set. We review this result below. To this end, we recall some concepts and results of Doignon and Falmagne (1999), which we adapt to the general case in which the empty set is not assumed to belong to the family[5]. Even though the proofs of Theorems 10 and 11 are essentially those of Theorems 1.25 and 1.26 in Doignon and Falmagne (1999), we include those proofs for completeness because our context is more general.

**9 Definition.** For any $x$ in $\mathcal{X} = \cup \mathcal{F}$, where $\mathcal{F}$ is a $\cup$-closed family, an *atom at $x$* is a minimal set of $\mathcal{F}$ containing $x$ (where 'minimal' is with respect to set inclusion). A set $X$ in $\mathcal{F}$ is called an *atom*[6] if either $X = \varnothing \in \mathcal{F}$, or there is some $x \in \mathcal{X}$ such that $X$ is an atom at $x$. Writing $\mathfrak{P}(\mathcal{F})$ for the power set of $\mathcal{F}$, we denote by $\text{ß}(x)$ the collection of all the atoms at $x$ and refer to $\text{ß} : \mathcal{X} \to \mathfrak{P}(\mathcal{F})$ as the *surmise function* of $\mathcal{F}$. Clearly, since $\mathcal{X}$ is finite, we have $\text{ß}(x) \neq \varnothing$ for every $x \in \mathcal{X}$; thus, there is at least one atom at every point of $\mathcal{X}$. (But see Counterexample 12.)

**10 Theorem.** *A nonempty set $X$ in a $\cup$-closed family $\mathcal{F}$ is an atom if and only if $X \in \mathcal{H}$ for any subfamily $\mathcal{H}$ of $\mathcal{F}$ satisfying $\cup \mathcal{H} = X$.*

---

[5] All of Doignon and Falmagne (1999)'s results were developed in the context of knowledge spaces, that is $\cup$-closed families containing the empty set. We drop the latter condition here.

[6] Our meaning of the term 'atom' is different from its usage in lattice theory; cf. Birkhoff (1967), Davey and Priestley (1990). It also slightly differs from that in Doignon and Falmagne (1999) because we do not allow the empty union of a family (see Footnote 4, 5 and 6).

PROOF. (Necessity.) Suppose that $X$ is an atom at some $x \in \cup\mathcal{F}$, with $X = \cup\mathcal{H}$ for some subfamily $\mathcal{H}$ of $\mathcal{F}$. Then $x \in Y$ for some $Y \in \mathcal{H}$, with necessarily $Y \subseteq X$. This implies $Y = X$ because $X$ is a minimal set containing $x$, and so $X \in \mathcal{H}$. The case of the atom $\varnothing$ is straightforward.

(Sufficiency.) If some $X \in \mathcal{F}$ is not an atom, then for each $x \in X$, we must have $x \in Y(x) \subset X$ for some $Y(x) \in \mathcal{F}$. Writing $\mathcal{H} = \{Y(x) \,|\, x \in X\}$, we get $\cup\mathcal{H} = X$, with $X \notin \mathcal{H}$. $\qquad\square$

**11 Theorem.** *The base of a $\cup$-closed family $\mathcal{F}$ is the collection of all its atoms.*

PROOF. Let $\mathcal{A}$ be the collection of all the atoms of $\mathcal{F}$. We claim that $\mathcal{A}$ must be the base of $\mathcal{F}$. If $\varnothing \in \mathcal{F}$, we have $\varnothing \in \mathcal{A}$ by definition with $\cup\{\varnothing\} = \varnothing$. Notice that, for any $X \neq \varnothing$ in $\mathcal{F}$, the set $\mathcal{A}_X = \{Y \in \mathcal{A} \,|\, \exists x \in X, \ x \in Y \subseteq X\}$ exists because there is an atom at every point of $X \subseteq \mathcal{X}$. We have thus $\cup\mathcal{A}_X = X$ and so $\mathcal{A}$ spans $\mathcal{F}$ (whether or not $\varnothing \in \mathcal{F}$). Let now $\mathcal{H}$ be another subfamily of $\mathcal{F}$ spanning $\mathcal{F}$. Take any $Z \in \mathcal{A}$. Since $\mathcal{H}$ spans $\mathcal{F}$, there must be a subfamily $\mathcal{G}$ of $\mathcal{H}$ such that $\cup\mathcal{G} = Z$. By Theorem 10, we must have $Z \in \mathcal{G} \subseteq \mathcal{H}$; this yields $\mathcal{A} \subseteq \mathcal{H}$. Thus, $\mathcal{A}$ is a minimal family spanning $\mathcal{F}$ and so is the (unique) base of $\mathcal{F}$. $\qquad\square$

Note in passing that, in the infinite case, there may not be an atom at every point of the ground set $\mathcal{X} = \cup\mathcal{F}$ of a $\cup$-closed family $\mathcal{F}$. We already gave the example of the collection of all the open sets of $\mathbb{R}$. Below is another, simple example.

**12 A Counterexample.** Consider the infinite family $\mathcal{F} = \mathcal{G} + \mathcal{H}$, with

(6) $$\mathcal{G} = \{G_n \,|\, G_n = \{\ldots, \frac{1}{n+1}, \frac{1}{n}\}, \ n > 1\}$$

(7) $$\mathcal{H} = \{H_n \,|\, H_n = G_n + \{1\}, \ G_n \in \mathcal{G}\}.$$

The family $\mathcal{F}$ is $\cup$-closed and well-graded and there is no atom at 1. It is easily verified that this $\cup$-closed family $\mathcal{F}$ has no base. The $\cup$-closed family $\mathcal{H}$ is its own base and has no atom at 1 either.

We turn to Koppen (1998)'s result, which is formulated as the last statement in the theorem below. We recall that $\check{\mathcal{B}} = \mathcal{B} \setminus \{\varnothing\}$ for the base $\mathcal{B}$ of a learning space.

**13 Theorem.** *Suppose that $\mathcal{F}$ is a learning space with base $\mathcal{B}$ and surmise function ß. Then $\{ß(x) \,|\, x \in \cup\mathcal{F}\}$ is a partition of $\check{\mathcal{B}}$ if and only if there is a tight path from $\varnothing$ to $K$ for any $K \in \check{\mathcal{B}}$. Accordingly, $\mathcal{F}$ is well-graded if and only if $\{ß(x) \,|\, x \in \cup\mathcal{F}\}$ is a partition of $\check{\mathcal{B}}$.*

PROOF. **Observation.** By Theorem 11, for any $K \in \check{\mathcal{B}}$, there is some $y \in K$ such that $K$ is an atom at $y$. Moreover, the hypothesis that $\{ß(x) \,|\, x \in \cup\mathcal{F}\}$ is a partition of $\check{\mathcal{B}}$ and $|K| > 1$ implies that there exists, for any $y' \in K$ distinct from $y$, at least one atom at $y'$ strictly included in $K$. (Otherwise, we would have $K \in ß(y) \cap ß(y')$.)

Assume that $\{ß(x) \,|\, x \in \cup\mathcal{F}\}$ is a partition of $\check{\mathcal{B}}$. Take any $K$ in $\check{\mathcal{B}}$ and suppose that $|K| = n$. By the Observation, $K_n = K$ is an atom at $x_n$ for some $x_n \in K$. We

6

use induction on $n$. If $n = 1$, then $\varnothing, \{x_1\} = K$ is the tight path. Suppose that we have a tight path $K_0 = \varnothing, K_1 = \{x_1\}, \ldots, K_j = \{x_1, \ldots, x_j\}$ from $\varnothing$ to $K_j \subset K_n = K$. From the Observation, we know that there exists an atom $L_\ell$ at $y_\ell$ for any $y_\ell \in K \setminus K_j$, with $1 \leq \ell \leq n - j$ and $L_\ell \subset K$. If $|K_j \cup L_\ell| > j + 1$ for some index $\ell$, then, again by the Observation, there is some index $i \neq \ell$, $1 \leq i \leq n - j$, such that $L_i \subset L_\ell$ is an atom at $y_i \in L_\ell \setminus K_j$, with $j + 1 \leq |K_j \cup L_i| < |K_j \cup L_\ell|$. By elimination, we have necessarily some $y_k \notin K_j$, $1 \leq k \leq n - j$ and an atom $L_k \subset K$ at $y_k$ such that $K_j \cup \{y_k\} = K_j \cup L_k$. Defining $x_{j+1} = y_k$ and $K_{j+1} = K_j \cup \{x_{j+1}\}$, we obtain the tight path $K_0 = \varnothing, K_1, \ldots, K_{j+1}$ from $\varnothing$ to $K_{j+1} \subseteq K$. Applying induction yields the necessity in the first statement.

Conversely, assume that there is a tight path from $\varnothing$ to $L$ for any $L \in \mathcal{B}$. Suppose that $\varnothing \neq K \in \mathcal{B}(x) \cap \mathcal{B}(y)$ for some $K \in \check{\mathcal{B}}$ and some distinct $x, y \in \cup \mathcal{F}$. A contradiction ensues because no tight path

$$K_0 = \varnothing, K_1 = \{x_1\}, \ldots, K_n = \{x_1, \ldots, x_n\} = K$$

from $\varnothing$ to $K$ can exist. Indeed, we must have $x, y \in K \setminus K_{n-1}$ since $K$ is an atom at both $x$ and $y$; and yet $|K| = n > n - 1 = |K_{n-1}|$.

The last statement of the theorem follows from the last statement in Theorem 7. $\square$

As announced, this result does not generalize to the case in which the family $\mathcal{F}$ does not contain the empty set, even if we assume that the family is *discriminative* that is, satisfies the condition: for all $x, y \in \cup \mathcal{F}$

$$(\forall X \in \mathcal{F}, \ x \in X \Leftrightarrow y \in X) \iff x = y.$$

**14 A Counterexample.** Consider the family $\mathcal{K}$ defined by the base

$$\mathcal{A} = \{\{x, y, c\}, \{y, d\}, \{c, d\}\}.$$

We get the surmise function

$$\mathcal{B}(x) = \{\{x, y, c\}\}, \quad \mathcal{B}(y) = \{\{x, y, c\}, \{y, d\}\},$$
$$\mathcal{B}(c) = \{\{x, y, c\}, \{c, d\}\}, \quad \mathcal{B}(d) = \{\{y, d\}, \{c, d\}\}.$$

It is easily checked that $\mathcal{K}$ is discriminative and well-graded; yet, the surmise function does not define a partition of the base $\mathcal{A}$.

# Algorithms

In the algorithms described in this section, we are given as input a family of sets $\mathcal{B}$, which is purported to be the base of a $\cup$-closed family $\mathcal{F}$. We wish to test whether this is true, and if so to determine other properties of $\mathcal{F}$ such as whether it is well-graded or a learning space. In many cases the definitions given in earlier sections of this paper may already be directly translated into algorithms, but a definition may be translated into an algorithm in multiple ways, some more efficient than others; the content of

the results lies less in the pure existence of the algorithms and more in designing the algorithms so they perform their tasks efficiently and in analyzing how much time they take to run. The time for our algorithms should be polynomial in the size of our input, if possible; this size is the sum of the cardinalities of the sets in $\mathcal{B}$. In particular, this requirement for polynomial time precludes explicit construction of $\mathcal{F}$ as the span of $\mathcal{B}$, as $\mathcal{F}$ may have exponentially greater size.

We assume a standard random-access-machine model of computation in which simple arithmetic steps and memory access operations may be performed in constant time. The input to our algorithms will be families of sets. We assume that each set element is represented as an object that takes a constant amount of computer storage and with which additional information may be associated. For instance, a natural representation with these properties would be to represent the $n$ elements of a set family as integers in the range from 0 to $n-1$; we may then associate information with each element by using these integers as array indices. We represent an input set as a list of elements, and an input set family as a list of lists of elements. As is standard in the analysis of algorithms, we use $O$-notation to simplify the stated time bounds for our algorithms.

**15 Definition.** In order to analyze and compare the running times of our algorithms, we need parameters to describe the input size. We define $n$ to be the number of sets in $\mathcal{B}$, $\ell$ to be the size of the largest set in $\mathcal{B}$, and $m$ to be the sum of cardinalities of sets in $\mathcal{B}$. We say that an algorithm runs in polynomial time if its worst-case running time can be upper bounded by a polynomial function of $\ell$, $m$ and $n$. For purposes of comparing run times it is convenient to note that $\ell \leq m \leq n\ell$.

**16 Definition.** The *endpoints* of a set $X$ belonging to a base $\mathcal{B}$ are the elements of the set

$$X \setminus \bigcup_{Y \in \mathcal{B}, Y \subset X} Y.$$

That is, the endpoints of $X$ are the elements of $X$ that are not contained in any proper subset of $X$ that belongs to $\mathcal{B}$. Equivalently, $x$ is an endpoint of a set $X$ in a base $\mathcal{B}$ if $X$ is an atom at $x$.

**17 Lemma.** *There is an algorithm that takes as input a set family $\mathcal{B}$ and a set $X \in \mathcal{B}$, and that outputs the endpoints of $X$, using time $O(m)$.*

PROOF. We associate with each element $x$ in $\cup\mathcal{B}$ a Boolean variable that is true if and only if $x$ is in $X$; setting up these variables takes time $O(m)$. By examining the value for $x$, we may test whether $x$ belongs to $X$ in constant time. For each set $Y \in \mathcal{B}$, we use these bits to determine whether $Y \subset X$, by testing each of the members of $Y$, in time $O(|\cup Y|)$. By performing this test for all sets in $\mathcal{B}$, we may determine a collection of the subsets of $X$ that are in $\mathcal{B}$, in total time $O(m)$. We then associate a second Boolean variable with each member of $X$; initially we set all of these variables to false. For each $Y$ in our collection of subsets of $X$, we loop through the elements of $Y$, and set the Boolean variables associated with each of these elements to true. Finally, we loop through the elements of $X$, and form a list of the elements for which the associated Boolean value remains false. These elements are the endpoints of $X$.

The runtime of this algorithm is dominated by the steps in which we find the subsets of $X$ and then use those subsets to mark covered elements of $X$; both of these steps take $O(m)$ total time. $\square$

**18 Theorem.** *Given a family $\mathcal{B}$ of sets, we may determine in time $O(nm)$ whether $\mathcal{B}$ is the base of a $\cup$-closed family $\mathcal{F}$.*

PROOF. We use Lemma 17 to calculate the endpoints of each $X \in \mathcal{B}$. By definition, $X$ is an atom if and only if it is empty or has a nonempty set of endpoints; thus, by Theorem 11, $\mathcal{B}$ is the base of its span if and only if every set in $\mathcal{B}$ is either empty or has a nonempty set of endpoints. There are $n$ sets, each of which takes time $O(m)$ to test, so the total time is $O(nm)$. $\square$

**19 Lemma.** *Suppose that set family $\mathcal{B}$ contains the empty set. Then $\mathcal{B}$ is the base of a $\cup$-closed well-graded family if and only if each nonempty $X \in \mathcal{B}$ has one endpoint.*

This is closely related to some results of Koppen (1998) (see also Doignon and Falmagne, 1999, Theorem 3.15, Condition (ii)).

PROOF. If some $X \in \mathcal{B}$ has two or more endpoints $x$ and $y$, then $X$ belongs to both $\sigma(x)$ and $\sigma(y)$, so the surmise function $\sigma$ is not a partition of $\check{\mathcal{B}}$. If some nonempty $X$ has no endpoint, it is not an atom and not part of a base. Conversely if every nonempty $X \in \mathcal{B}$ has one endpoint, then $\sigma$ partitions $\check{\mathcal{B}}$ according to those endpoints. The result follows from Theorem 13. $\square$

**20 Theorem.** *Given a family $\mathcal{B}$ of sets, we may determine in time $O(nm)$ whether $\mathcal{B}$ is the base of a learning space.*

PROOF. We first check that $\mathcal{B}$ contains the empty set; if not, it cannot be the base of a learning space. Then, as in Theorem 18, we apply Lemma 17 to calculate the endpoints of each $X \in \mathcal{B}$. By Lemma 19, $\mathcal{B}$ is the base of a $\cup$-closed well-graded family if and only if each nonempty $X \in \mathcal{B}$ has exactly one endpoint. There are $n$ sets, each of which takes time $O(m)$ to test, so the total time is $O(nm)$. $\square$

**21 Definition.** For any set family $\mathcal{B}$ and any set $X \in \mathcal{B}$, let $\mathcal{B}/X$ denote the family of sets $\{Y \setminus X \,|\, Y \in \mathcal{B}\}$.

**22 Lemma.** *Let $\mathcal{B}$ be the base of a $\cup$-closed family $\mathcal{F}$. Then $\mathcal{F}$ is well-graded if and only if, for each $X$ in $\mathcal{B}$, the family $\mathcal{B}/X$ spans a learning space.*

PROOF. A tight path in $\mathcal{F}$ from $X$ to some set $Y \supset X$ corresponds (via set-theoretic difference of each path member with $X$) to a tight path in $\mathcal{F}/X$ from the empty set to $Y \setminus X$. Conversely, a tight path in $\mathcal{F}/X$ from the empty set to $Y \setminus X$ corresponds (via set-theoretic union of each path member with $X$) to a tight path in $\mathcal{F}$ from $X$ to $Y$. The result follows from Theorem 7. $\square$

**23 Theorem.** *Given a family $\mathcal{B}$ of sets, we may determine in time $O(n^2m)$ whether $\mathcal{B}$ is the base of a $\cup$-closed well-graded family.*

PROOF. We may first test whether $\mathcal{B}$ is a base by Theorem 18. Next, for each $X \in \mathcal{B}$, we form the set $\mathcal{B}_X$ consisting of the empty set and the sets in $\mathcal{B}/X$ that have a nonempty set of endpoints with respect to $\mathcal{B}/X$, and test whether $\mathcal{B}_X$ is the base of a learning space by Theorem 20. $\mathcal{B}$ itself is the base of a $\cup$-closed well-graded family if and only if each $\mathcal{B}_X$ passes this test, by Lemma 22. There are $n$ sets $\mathcal{B}_X$, each takes time $O(nm)$ to construct and test, and so the total time bound is $O(n^2m)$. $\square$

We now consider the situation in which $\mathcal{B}$ is not itself the base of a well graded family. Can we modify $\mathcal{B}$ to produce a well graded family that is as close as possible, in some sense, to the span of $\mathcal{B}$?

**24 Definition.** A *minimal well-graded extension* of a family of sets $\mathcal{B}$ is a well-graded $\cup$-closed set family $\mathcal{F}$ such that $\mathcal{B} \subset \mathcal{F}$, and such that no $\cup$-closed $\mathcal{F}'$ with $\mathcal{B} \subset \mathcal{F}' \subset \mathcal{F}$ is well-graded. A *path family* for a family of sets $\mathcal{B}$ is a set paths $\pi_{K,L}$ with $K$ and $L$ in $\mathcal{B}$; $\pi_{K,L}$ may use sets not belonging to the span of $\mathcal{B}$, but is required to be a tight path in the power set of $\cup\mathcal{B}$. We observe that the length of a path $\pi_{K,L}$ is at most the cardinality of $L$, and therefore that the total length of all paths in a path family is $O(nm)$. A *path extension* $\mathcal{F}$ of $\mathcal{B}$ is formed from a path family by letting $\mathcal{B}'$ consist of all the sets occurring on paths $\pi_{K,L}$ and letting $\mathcal{F}$ be the span of $\mathcal{B}'$.

**25 Lemma.** *Any path extension is well-graded.*

PROOF. We show that, for every $K$ and $L$ in $\mathcal{B}'$, where $\mathcal{B}'$ is the family of sets occurring on paths $\pi_{X,Y}$ in a path extension of $\mathcal{B}$, that there exists a tight path in the span of $\mathcal{B}'$ from $K$ to $K \cup L$.

Thus, suppose $K$ belongs to a path $\pi_{A,B}$ and $L$ belongs to a path $\pi_{C,D}$. To form a tight path from $K$ to $K \cup L$ in $\mathcal{F}$, we concatenate the following three paths:

1. a tight path from $A$ to $K$ along path $\pi_{A,B}$,

2. a tight path from $K$ to $K \cup C$, formed by the union of $K$ with the sets in path $\pi_{A,C}$, and

3. a tight path from $K \cup C$ to $K \cup L$, formed by the union of $K \cup C$ with the sets in the portion of path $\pi_{C,D}$ that extends from $C$ to $L$.

When this concatenation would cause the same set to appear repeatedly, we discard the duplicate sets. It is straightforward to verify that each set in this concatenation of paths belongs to the span of $\mathcal{B}'$. Thus, we can form a tight path from any $K$ to $K \cup L$ in this span, and therefore, by Theorem 7, the span is well-graded. $\square$

**26 Lemma.** *Any minimal well-graded extension is a path extension.*

PROOF. Let $\mathcal{F}$ be a minimal well-graded extension of $\mathcal{B}$. Then by Theorem 7 we can find a path family for $\mathcal{B}$, such that each set occurring in each path belongs to $\mathcal{F}$. By Lemma 25, the corresponding path extension is well-graded, and it is a subfamily of $\mathcal{F}$ and contains every set in $\mathcal{B}$. By the minimality of $\mathcal{F}$, this path extension must coincide with $\mathcal{F}$. $\square$

10

It is straightforward to combine the results above in an algorithm that finds a minimal well-graded extension of any set family in polynomial time: construct a path family arbitrarily, and then for each set in its base, determine whether the set can be removed by testing the result of the removal for well-gradedness, using Theorem 23 to do these tests. However the polynomial time bound of this algorithm would be large. We now describe a more efficient algorithm for the same task, based on a more careful choice of path family.

**27 Theorem.** *Given any family of sets $\mathcal{B}$, we can find a minimal well-graded extension of $\mathcal{B}$ in time $O(nm\ell + n^3m)$.*

PROOF. We simultaneously form the paths $\pi_{K,L}$ in a path family, and a superset $\mathcal{B}'$ of $\mathcal{B}$ that includes a base for our eventual well-graded extension, by adding sets in order by the cardinality of the sets. At step $i$ of the process, we include sets of cardinality $i$ into $\mathcal{B}'$, taking care as we do that all the sets we add are necessary for well-gradedness. As we do so, we maintain the following data:

- $S_{K,L}$ is the union of all sets $X \in \mathcal{B}'$ such that $X \subset K \cup L$ and $X \cup K \neq K \cup L$.

- $c_{A,B,C,D}$ is a Boolean value, true if and only if $S_{A,B} \subset C \cup D$.

In the $i$th step of the algorithm, we consider the set $\Pi_i$ of all paths $\pi_{K,L}$ such that $|S_{K,L}| = i - 1$ and such that $|K \cup L| > i$. We will add sets of cardinality $i$ to $\mathcal{B}'$, of the form $S_{K,L} \cup \{x\}$ for some $x$ in $(K \cup L) \setminus S_{K,L}$, in order to allow one more step on each path. However, note that if $S_{A,B} = S_{C,D}$ then a single set of this type may allow an additional step for multiple paths.

We observe that, if $\pi_{A,B}$ and $\pi_{C,D}$ are both in $\Pi_i$, then $c_{A,B,C,D}$ is true if and only if $S_{A,B} = S_{C,D}$. Thus the relation $c$ can be viewed as an equivalence relation on the paths in $\Pi_i$. As part of our calculation in step $i$ of the algorithm, we construct the equivalence classes of this equivalence relation.

For each equivalence class, we form a bipartite graph $(U, V, E)$. Here $U$ consists of pairs $(K, L)$ corresponding to paths $\pi_{K,L}$ in the equivalence class. $V$ consists of elements $x$ in the sets $(K \cup L) \setminus S_{K,L}$, for paths $\pi_{K,L}$ in the equivalence class. We draw an edge from $(K, L)$ to $x$ if $x \in (K \cup L) \setminus S_{K,L}$. We find a minimal subset of $V$ that dominates every vertex of $U$ in this graph; this gives us a minimal family of sets that we can add to $\mathcal{B}'$ in order to take another step on each path in the equivalence class. This minimal dominating set can be found by repeatedly either including in it a vertex in $V$ that is the only neighbor of some vertex in $U$ or, if no such vertex in $U$ exists, removing from the graph an arbitrarily chosen vertex in $V$; the total time to perform this step is proportional to the size of the graph.

Once we have found these sets to add to $\mathcal{B}'$, we must update the data we are maintaining so that we may repeat this computation for a larger value of $i$. Whenever we add a set corresponding to a path $\pi_{A,B}$ in $\Pi_i$ and element $x$, we examine all paths $\pi_{C,D}$ for which $c_{A,B,C,D}$ is true. If $x \in C \cup D$, we include $x$ as a new member of $S_{C,D}$. (In particular, $c_{A,B,A,B}$ will always be true, and we will always include $x$ as a new member of $S_{A,B}$.) However, if $x \notin C \cup D$, we instead set $c_{A,B,C,D}$ to false.

We now analyze the running time of this algorithm:

- We may compute the initial value of each set $S_{K,L}$ in time $O(m)$, simply by testing each other set $X$ in time $O(|X|)$, after an initial $O(|K| + |L|)$ time preprocessing stage to construct data structures for testing membership in $K \cup L$. and $L \setminus K$. Thus, we may construct all such sets in time $O(n^2 m)$.

- We may compute the initial value of $c_{A,B,C,D}$ in time $O(|A| + |B| + |C| + |D|)$. Adding this up over all quadruples $A, B, C, D$ produces a runtime of $O(n^3 m)$.

- Identifying $\Pi_i$ takes time $O(n^2)$. There are $O(m)$ steps of the algorithm, so the total time for this identification is $O(n^2 m)$.

- The sum of the cardinalities of the sets $\Pi_i$, summed over all $i$, is $O(nm)$, because each time we include a set in $\Pi_i$ we take a step on the corresponding path, and the total length of all paths in a path family is $O(nm)$. We may identify the equivalence class of a single path in time $O(n^2)$; therefore, the total time to construct equivalence classes, throughout the course of the algorithm, is $O(n^3 m)$.

- Each vertex in $U$ in the bipartite graph constructed for an equivalence class may have $O(\ell)$ neighbors. Therefore, the total size of all bipartite graphs so constructed, and the total time to find dominating sets in these graphs, is $O(nm\ell)$.

- Each set added to $\mathcal{B}'$ can be constructed explicitly, as a list of elements, from the data structures we already have, in time $O(\ell)$. Thus, the total time to list all these sets is $O(nm\ell)$. In addition, for each such set, we spend $O(n^2)$ time examining the paths for which $c_{A,B,C,D}$ is true, the total for which over the course of the algorithm is $O(n^3 m)$.

Thus, the total time for all of these steps is $O(nm\ell + n^3 m)$. $\qquad\square$

It may be seen as a flaw in the completion algorithm described above that not every set in the input family $\mathcal{B}$ is necessarily part of a base of the output family $\mathcal{F}$ it produces. Given $\mathcal{B}$, can we find a wg-family $\mathcal{F}$ such that every set in $\mathcal{B}$ is part of the base of $\mathcal{F}$? Unfortunately, as we now show, this problem appears to be intractable.

**28 Theorem.** *It is NP-complete, given a set family $\mathcal{B}$, to determine whether there exists a well-graded $\cup$-closed set family $\mathcal{F}$ such that $\mathcal{B}$ is a subset of the base of $\mathcal{F}$.*

PROOF. If an $\mathcal{F}$ satisfying this requirement exists, we can choose $\mathcal{F}$ to be minimal and therefore, by Lemma 26, a path extension. Thus, we can test in NP whether $\mathcal{F}$ exists by nondeterministically choosing a path extension, applying Lemma 17 to find the endpoints of all sets included on paths in the extension, and verifying that each member of $\mathcal{B}$ has an endpoint. Therefore, determining whether $\mathcal{F}$ exists belongs to NP, the easier part of proving that it is NP-complete.

To finish the NP-completeness proof, we reduce the problem from a known NP-complete problem, 3-satisfiability (Garey and Johnson, 1979). A 3-satisfiability instance consists of sets of *variables* $V$, complements of variables $\bar{V} = \{\bar{v} \mid v \in V\}$, and a set $C$ of clauses, where each clause is a set of three terms, and where a term is any element of $V \cup \bar{V}$. A truth assignment is any function $f$ from $V$ to $\{0, 1\}$; we may

extend $f$ to the domain $V \cup \bar{V}$ by $f(\bar{v}) = 1 - f(v)$. A truth assignment is *satisfying* if each clause of $C$ contains at least one term mapped by $f$ to 1, and a 3-satisfiability instance is *satisfiable* if and only if it has a satisfying assignment.

From a 3-satisfiability instance $(V, \bar{V}, C)$ we form a set family $\mathcal{B}$, the ground set of which will be the terms and clauses of the instance: $\cup \mathcal{B} = V \cup \bar{V} \cup C$. For each variable $v \in V$, we include in $\mathcal{B}$ the set $\{v, \bar{v}\}$, and for each clause $c$ corresponding to the conjunction of three terms $u$, $v$, and $w$ we include in $\mathcal{B}$ two sets, $\{c\}$ and $\{c, u, v, w\}$. Additionally, we include in $\mathcal{B}$ the empty set. As we now show, the resulting set $\mathcal{B}$ forms a subset of the base of a well-graded $\cup$-closed set family $\mathcal{F}$, if and only if the given 3-satisfiability instance is satisfiable.

In one direction, suppose we have a satisfying truth assignment $f$. Let $\mathcal{H} = \{\{x\} \mid x \in V \cup \bar{V} \text{ and } f(x) = 0\}$. Let $t_i(c)$, for $i \in \{0, 1, 2\}$ map clauses to terms in such a way that $c = \{t_0(c), t_1(c), t_2(c)\}$ and $f(t_2(c)) = 1$. Let $T_0 = \{\{c, t_0(c)\} \mid c \in C\}$ and $T_1 = \{\{c, t_0(c), t_1(c)\} \mid c \in C\}$. We let $\mathcal{F}$ be the span of $\mathcal{H} \cup \mathcal{B} \cup T_0 \cup T_1$. For each set $\{v, \bar{v}\} \in \mathcal{B}$ in there is a tight path in $\mathcal{F}$ from the empty set through $\{v\}$ or $\{\bar{v}\}$ respectively as $v$ or $\bar{v}$ is mapped by $f$ to 0; the element of $\{v, \bar{v}\}$ not mapped to 0 is an endpoint of $\{v, \bar{v}\}$. For each set $\{c, t_0(c), t_1(c), t_2(c)\} \in \mathcal{B}$ there is a tight path in $\mathcal{F}$ from the empty set through sets $\{c\}, \{c, t_0(c)\}, \{c, t_0(c), t_1(c)\}$, and $t_2(c)$ is an endpoint. The paths through $\{c\}, \{c, t_0(c)\}, \{c, t_0(c), t_1(c)\}$ also form tight paths in $\mathcal{F}$ to each set in $T_0$ and $T_1$. Thus, every set in $\mathcal{H} \cup \mathcal{B} \cup T_0 \cup T_1$ has a tight path in $\mathcal{F}$ from the empty set, so by Theorem 7 $\mathcal{F}$ is well-graded, and every nonempty set in $\mathcal{B}$ has an endpoint, so by Lemma 19 $\mathcal{B}$ is part of the base of $\mathcal{F}$. Thus, we have shown that, if $f$ is a satisfying truth assignment, $\mathcal{B}$ forms a subset of the base of a well-graded $\cup$-closed family.

In the other direction, suppose that there exists a set family $\mathcal{F}$ that is a minimal path extension of $\mathcal{B}$ for which $\mathcal{B}$ is a subset of the base. Then, in order to have a tight path from the empty set to $\{v, \bar{v}\}$, while not eliminating that set from the base, $\mathcal{F}$ must contain exactly one of the two sets $\{v\}$, $\{\bar{v}\}$; form a truth assignment $f$ in which we assign $v$ the value 1 if $\{\bar{v}\}$ is in $\mathcal{F}$ and the value 0 if $\{v\}$ is in $\mathcal{F}$. This must be a satisfying assignment, for if a clause $c$ had no variable satisfying it then the set $\{c, u, v, w\}$ would have no endpoints and therefore couldn't be part of the base of $\mathcal{F}$. Thus, we have shown that, if $\mathcal{B}$ forms a subset of the base of a well-graded $\cup$-closed family, then the 3-satisfiability instance $(V, \bar{V}, C)$ has a satisfying assignment.

We have described a polynomial time many-one reduction from the known NP-complete problem of 3-satisfiability to the problem of testing whether a set family is a subset of a base of a well-graded $\cup$-closed family, and we have shown that the latter problem is in NP. Therefore, it is NP-complete. $\square$

**29 Remark.** Since the family $\mathcal{B}$ formed by this reduction contains the empty set, the same reduction shows that it is also NP-complete to determine whether a given set $\mathcal{B}$ is a subset of the base of a learning space.

**30 Remark.** It is natural to desire, not just a minimal well-graded extension, but an extension that is minimum, either in the sense of having the smallest cardinality as a well-graded set family, in the sense of having the smallest cardinality base, in the sense

of having the smallest number of additional sets added to the input family $\mathcal{B}$, or in the sense of minimizing the sum of cardinalities of additional sets or of the base. We expect that these problems are computationally intractable, but do not have a hardness proof for them. The problems of minimizing base size at least belong to NP, but minimizing the cardinality of $\mathcal{F}$ may not since it involves counting the members of a family of sets that may be exponentially larger than the input.

# References

G. Birkhoff. *Lattice Theory*. American Mathematical Society, Providence, R.I., 1967.

K.P. Bogart. Maximal dimensional partially ordered sets. I. Hiraguchi's theorem. *Discrete Mathematics*, 5:21–31, 1973.

E. Cosyn and H.B. Uzun. Axioms for learning spaces. In press in the *Journal of Mathematical Psychology*, 2008.

B.A. Davey and H.A. Priestley. *Introduction to Lattices and Order*. Cambridge University Press, Cambridge, London, and New Haven, 1990.

D.Z. Djoković. Distance preserving subgraphs of hypercubes. *Journal of Combinatorial Theory, Ser. B*, 14:263–267, 1973.

J.-P. Doignon and J.-Cl. Falmagne. Well-graded families of relations. *Discrete Mathematics*, 173:35–44, 1997.

J.-P. Doignon and J.-Cl. Falmagne. *Knowledge Spaces*. Springer-Verlag, Berlin, Heidelberg, and New York, 1999.

J.-P. Doignon and J.-Cl. Falmagne. Spaces for the Assessment of Knowledge. *International Journal of Man-Machine Studies*, 23:175–196, 1985.

D. Eppstein, J.-Cl. Falmagne, and Ovchinnikov S. *Media Theory*. Springer-Verlag, Berlin, Heidelberg, and New York, 2007.

J.-Cl. Falmagne. Stochastic token theory. *Journal of Mathematical Psychology*, 41(2): 129–143, 1997.

J.-Cl. Falmagne and J.-P. Doignon. Stochastic evolution of rationality. *Theory and Decision*, 43:107–138, 1997.

J.-Cl. Falmagne, E. Cosyn, J.-P. Doignon, and N. Thiéry. The assessment of knowledge, in theory and in practice. In B. Ganter and L. Kwuida, editors, *Formal Concept Analysis, 4th International Conference, ICFCA 2006, Dresden, Germany, February 13–17, 2006*, Lecture Notes in Artificial Intelligence, pages 61–79. Springer-Verlag, Berlin, Heidelberg, and New York, 2006.

M.R. Garey and D.S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freemann, 1979.

R.L. Graham and H. Pollak. On addressing problem for loop switching. *Bell Systems Technical Journal*, 50:2495–2519, 1971.

W. Imrich and S. Klavžar. *Product Graphs*. John Wiley & Sons, London and New York, 2000.

M. Koppen. On alternative representations for knowlede spaces. *Mathematical Social Sciences*, 36:127–143, 1998.

B. Korte, L. Lovász, and R. Schrader. *Greedoids*. Number 4 in Algorithms and Combinatorics. Springer-Verlag, 1991.

V.B. Kuzmin and S. Ovchinnikov. Geometry of preference spaces I. *Automation and Remote Control*, 36:2059–2063, 1975.

S. Ovchinnikov. Convexity in subsets of lattices. *Stochastica*, IV:129–140, 1980.

P.M. Winkler. Isometric embedding in products of complete graphs. *Discrete Applied Mathematics*, 7:221–225, 1984.